

# Conditioning of implicit Runge-Kutta integration for finite element approximation of linear diffusion equations on anisotropic meshes\*

Weizhang Huang<sup>†</sup>    Lennard Kamenski<sup>‡</sup>    Jens Lang<sup>§</sup>

March 19, 2017

## Abstract

The conditioning of implicit Runge-Kutta (RK) integration for linear finite element approximation of diffusion equations on general anisotropic meshes is investigated. Bounds are established for the condition number of the resulting linear system with and without diagonal preconditioning for the implicit Euler (the simplest implicit RK method) and general implicit RK methods. It is shown that the conditioning of an implicit RK method behaves like that of the implicit Euler method. Two solution strategies are considered for the linear system resulting from general implicit RK integration: the simultaneous solution where the system is solved as a whole and a successive solution which follows the commonly used implementation of implicit RK methods to first transform the system into a number of smaller systems using the Jordan normal form of the RK matrix and then solve them successively. The obtained bounds for the condition number have explicit geometric interpretations and take the interplay between the diffusion matrix and the mesh geometry into full consideration. They show that there are three mesh-dependent factors that can affect the conditioning: the number of elements, the mesh nonuniformity measured in the Euclidean metric, and the mesh nonuniformity with respect to the inverse of the diffusion matrix. They also reveal that the preconditioning using the diagonal of the system matrix, the mass matrix, or the lumped mass matrix can effectively eliminate the effects of the mesh nonuniformity measured in the Euclidean metric. Illustrative numerical examples are given.

**Keywords:** finite element method, anisotropic mesh, condition number, parabolic equation, implicit Runge-Kutta method

**AMS 2010 MSC:** 65M60, 65M50, 65F35, 65F15

---

\*Supported in part by the DFG under grant KA 3215/2-1 and the Darmstadt Graduate Schools of Excellence *Computational Engineering and Energy Science and Engineering*.

<sup>†</sup>Department of Mathematics, The University of Kansas, Lawrence, KS 66045, USA (whuang@ku.edu).

<sup>‡</sup>(kamenski@wias-berlin.de).

<sup>§</sup>Department of Mathematics, TU Darmstadt, 64293 Darmstadt, Germany (lang@mathematik.tu-darmstadt.de).

# 1 Introduction

The nonuniformity of adaptive meshes has considerable effects on the conditioning of the discrete approximation of partial differential equations (PDEs) and their efficient solution. To study these effects, we investigate the implicit Runge-Kutta (RK) integration for the linear finite element (FE) approximation of linear diffusion equations on general simplicial anisotropic meshes for the initial-boundary value problem (IBVP)

$$\begin{cases} \partial_t u = \nabla \cdot (\mathbb{D} \nabla u), & \mathbf{x} \in \Omega, \quad t > 0, \\ u(\mathbf{x}, t) = 0, & \mathbf{x} \in \partial\Omega, \quad t > 0, \\ u(\mathbf{x}, 0) = u_0(\mathbf{x}), & \mathbf{x} \in \Omega, \end{cases} \quad (1)$$

where  $\Omega \subset \mathbb{R}^d$  ( $d \geq 1$ ) is a bounded polygonal or polyhedral domain,  $u_0$  is a given initial solution, and  $\mathbb{D}$  is the diffusion matrix. We assume that  $\mathbb{D} = \mathbb{D}(\mathbf{x})$  is time independent and symmetric and uniformly positive definite on  $\Omega$ , i.e.,

$$\exists d_{\min}, d_{\max} > 0: \quad d_{\min} I \leq \mathbb{D}(\mathbf{x}) \leq d_{\max} I, \quad \forall \mathbf{x} \in \Omega, \quad (2)$$

where the less-than-or-equal sign means that the difference between the two matrices is positive semidefinite. We consider the Dirichlet boundary condition in this work but the analysis is applicable to other boundary condition types without major modifications.

Much effort has been made in the past to understand the effects of mesh nonuniformity on the conditioning of FE approximations. For example, Fried [7] obtains a bound on the condition number of the stiffness matrix for the linear FE approximation of the Laplace operator for general meshes. For the Laplace operator on isotropic adaptive grids, Bank and Scott [2] show that the condition number of the diagonally scaled stiffness matrix is essentially the same as for a regular mesh. Ainsworth, McLean, and Tran [1] and Graham and McLean [8] extend this result to the boundary element equations for locally quasi-uniform meshes and provide a bound in terms of patch volumes and aspect ratios. Du et al. [5] obtain a bound on the condition number of the stiffness matrix for a general diffusion operator on anisotropic meshes which reveals the relation between the condition number and some mesh quality measures. For the FE approximations of parabolic problems, Zhu and Du [20, 21] develop mesh dependent stability and condition number estimates for the explicit and implicit Euler methods. In the case of a lumped mass matrix ( $M_{lump}$ ), Shewchuk [19] provides a bound on the largest eigenvalue of  $M_{lump}^{-1} A$  in terms of the maximum eigenvalues of local element matrices, where  $A$  is the stiffness matrix. The results mentioned above allow anisotropic adaptive meshes but do not fully take into account the interplay between the mesh geometry and the diffusion matrix. (An exception is the bound by Shewchuk which includes the effects of the diffusion coefficients.) Moreover, the existing analysis employs mesh restrictions in form of mesh regularity assumptions, e.g., local mesh uniformity, or parameters in final estimates that are related to mesh regularity, such as the maximum ratio of volumes of neighboring elements or the maximum number of elements in a patch.

The objective of this work is to develop estimates on the conditioning of the resulting linear system that take the interplay between the mesh geometry and the diffusion matrix into full consideration, have explicit geometric interpretation, and make no prior assumptions on the mesh regularity. This is a continuation of our previous effort to develop bounds for the condition number of the stiffness matrix for the linear FE equations of a general diffusion operator on arbitrary anisotropic meshes [16, 17] and the largest permissible time steps for explicit RK schemes for both linear and high order FE approximations of the IBVP (1) [14, 15]. In particular, these bounds show [17] that the condition number of the stiffness matrix depends on three factors: the factor

depending on the number of mesh elements and corresponding to the condition number of the linear FE equations for the Laplace operator on a uniform mesh, the nonuniformity of the mesh viewed in the metric defined by the inverse diffusion matrix,  $\mathbb{D}^{-1}$ , and the mesh nonuniformity measured in the Euclidean metric. Moreover, the Jacobi preconditioning, an optimal diagonal scaling for a symmetric positive definite sparse matrix [11, Corollary 7.6ff.], can effectively eliminate the effects of mesh nonuniformity and reduce those of the mesh nonuniformity with respect to  $\mathbb{D}^{-1}$  [17]. Detailed characterizations of the condition number according to the mesh concentration distribution can be obtained using Green's functions [7, 16].

We first consider the implicit Euler method (the simplest implicit RK method) and establish bounds for the condition number of the corresponding system matrix with and without diagonal scaling. For general implicit RK methods, we consider two strategies for solving the resulting system: the simultaneous solution and a successive solution. For the simultaneous solution, the system is solved as a whole. For the successive solution, which follows the commonly used implementation of implicit RK methods [3, 4], the system is first transformed into a number of smaller systems using the Jordan normal form of the RK matrix and then solved successively. We show that the conditioning of the implicit RK integration is determined by the conditioning of two types of matrices. The first one is similar to the implicit Euler method and corresponds to the real eigenvalues of the RK matrix, while the second, twice as large, corresponds to the complex eigenvalues of the RK matrix (cf. (49) in Sect. 5). Obtained estimates reveal that the conditioning of implicit RK methods behaves like that of the implicit Euler method, both for the simultaneous and successive solution of the system.

The paper is organized as follows. We first introduce the FE formulation and its implicit RK integration (Sect. 2) and provide preliminary estimates for the extremal eigenvalues of the mass and stiffness matrices (Sect. 3). The main results for the conditioning of the coefficient matrices are given in Sects. 4 and 5, followed by numerical examples (Sect. 6) and conclusions (Sect. 7).

## 2 Linear FE approximation and implicit RK integration

Let  $\{\mathcal{T}_h\}$  be a given family of simplicial meshes for the domain  $\Omega$  and  $N$ ,  $N_v$ , and  $N_{vi}$  the number of mesh elements, vertices, and interior vertices, respectively. For convenience, we assume that the vertices are ordered such that the first  $N_{vi}$  vertices are the interior ones. The element patch associated with the  $j$ -th vertex is denoted by  $\omega_j$ ,  $K$  denotes a given mesh element and  $\hat{K}$  is the reference element which is assumed to have been taken as a unitary equilateral simplex. Element and patch volumes are denoted by  $|K|$  and  $|\omega_j| = \sum_{K \in \omega_j} |K|$ . For each mesh element  $K \in \mathcal{T}_h$ , we denote the invertible affine mapping from  $\hat{K}$  to  $K$  and its Jacobian matrix by  $F_K$  and  $F'_K$ , respectively. Note that  $F'_K$  is a constant matrix and  $\det(F'_K) = |K|$ .

Let  $V^h \subset H_0^1(\Omega)$  be the linear FE space associated with  $\mathcal{T}_h$ . The piecewise linear FE solution  $u_h(t) \in V^h$ ,  $t > 0$  for (1) is defined by

$$\int_{\Omega} \frac{\partial u_h}{\partial t} v_h \, d\mathbf{x} = - \int_{\Omega} (\nabla v_h)^T \mathbb{D} \nabla u_h \, d\mathbf{x}, \quad \forall v_h \in V^h, \quad t > 0, \quad (3)$$

subject to the initial condition

$$\int_{\Omega} u_h(\mathbf{x}, 0) v_h \, d\mathbf{x} = \int_{\Omega} u^0(\mathbf{x}) v_h \, d\mathbf{x}, \quad \forall v_h \in V^h.$$

It can be expressed as

$$u_h(\mathbf{x}, t) = \sum_{j=1}^{N_{vi}} u_j(t) \phi_j(\mathbf{x}),$$

where  $\phi_j$  is the linear basis function associated with the  $j$ -th vertex. Inserting this into (3) and taking  $v_h = \phi_k$ ,  $k = 1, \dots, N_{vi}$ , successively yields the system

$$M \frac{d\mathbf{u}}{dt} = -A\mathbf{u}, \quad (4)$$

where  $\mathbf{u} = (u_1, \dots, u_{N_{vi}})^T$  and  $M$  and  $A$  are the mass and stiffness matrices with

$$M_{kj} = \int_{\Omega} \phi_k \phi_j d\mathbf{x} \quad \text{and} \quad A_{kj} = \int_{\Omega} \nabla \phi_k \cdot \mathbb{D} \nabla \phi_j d\mathbf{x}, \quad k, j = 1, \dots, N_{vi}. \quad (5)$$

For the time integration of (4) we consider a general implicit  $s$ -stage RK method with the Butcher tableau

$$\begin{array}{c|cccc} c_1 & \gamma_{11} & \gamma_{12} & \cdots & \gamma_{1s} \\ c_2 & \gamma_{21} & \gamma_{22} & \cdots & \gamma_{2s} \\ \vdots & \vdots & \vdots & & \vdots \\ c_s & \gamma_{s1} & \gamma_{s2} & \cdots & \gamma_{ss} \\ \hline & b_1 & b_2 & \cdots & b_s \end{array}$$

and assume that the eigenvalues of the RK matrix have nonnegative real parts. This requirement is satisfied by most implicit RK methods (e.g., see Table 1). Applying the method to (4) yields

$$M\mathbf{v}_k + \Delta t A \sum_{j=1}^s \gamma_{kj} \mathbf{v}_j = -A\mathbf{u}^n, \quad k = 1, \dots, s, \quad (6)$$

$$\mathbf{u}^{n+1} = \mathbf{u}^n + \Delta t \sum_{k=1}^s b_k \mathbf{v}_k, \quad (7)$$

where  $\mathbf{u}^n$  and  $\mathbf{u}^{n+1}$  are the approximations of  $\mathbf{u}(t_n)$  and  $\mathbf{u}(t_{n+1})$ . The major cost of finding  $\mathbf{u}^{n+1}$  in the above method is the solution of (6) for  $\mathbf{v}_1, \dots, \mathbf{v}_s$ . Using the Kronecker matrix product  $\otimes$  (e.g., see [18]), the  $s \times s$  identity matrix  $I_s$ , and  $\Gamma := (\gamma_{kj})_{k,j=1}^s$ , the coefficient matrix of (6) can be expressed as

$$I_s \otimes M + \Delta t \Gamma \otimes A. \quad (8)$$

The simplest implicit RK method is the implicit Euler method with  $s = 1$ ,  $\Gamma = 1$ ,  $c_1 = 1$ , and  $b_1 = 1$  for which (6) and (7) are reduced to

$$(M + \Delta t A)\mathbf{u}^{n+1} = M\mathbf{u}^n. \quad (9)$$

In Sect. 4 we study the conditioning of the coefficient matrix  $M + \Delta t A$  related to the efficient iterative solution of (9). The system (9) can also be solved using a symmetric and positive definite preconditioner  $P$ , which leads to

$$P^{-\frac{1}{2}}(M + \Delta t A)P^{-\frac{1}{2}}\mathbf{w}^{n+1} = P^{-\frac{1}{2}}MP^{-\frac{1}{2}}\mathbf{w}^n, \quad \mathbf{w}^n = P^{\frac{1}{2}}\mathbf{u}^n. \quad (10)$$

The efficient iterative solution of (10) is related to the conditioning of  $P^{-\frac{1}{2}}(M + \Delta t A)P^{-\frac{1}{2}}$ . In this work, we consider as preconditioners the mass matrix  $M$ , its diagonal part  $M_D$ , the lumped mass matrix  $M_{lump}$ , and the diagonal part  $M_D + \Delta t A_D$  of  $M + \Delta t A$ . The choice  $P = M$  is of theoretical importance but impractical since  $M^{-\frac{1}{2}}$  is expensive to compute. The other choices are simple and economic to implement.

The results for the Euler method will be used in Sect. 5 to study the conditioning of (8) and its diagonally preconditioned version for general implicit RK methods.

### 3 Preliminary estimates on the extreme eigenvalues of the mass and stiffness matrices

Hereafter,  $C$  denotes a generic constant which may have different values at different appearances and may depend on the dimension, the choice of the reference element, and the reference basis linear functions but is independent of the mesh and the IBVP coefficients. For notation simplicity, when using this generic constant  $C$ , we will sometime write  $a \gtrsim b$  and  $a \lesssim b$  meaning  $a \geq Cb$  and  $a \leq Cb$ , respectively.

**Lemma 3.1** ([17, proof of Theorem 3.1]). *The mass matrix  $M$  and its diagonal part  $M_D$  are related by*

$$\frac{1}{2}M_D \leq M \leq \left(1 + \frac{d}{2}\right) M_D \quad \text{and} \quad M_{jj} = \frac{2|\omega_j|}{(d+1)(d+2)}. \quad (11)$$

**Lemma 3.2** ([15, Lemma 2.3]). *The mass matrix  $M$  and the lumped mass matrix  $M_{lump}$  are related by*

$$\frac{1}{d+2}M_{lump} \leq M \leq \frac{d+2}{2}M_{lump}. \quad (12)$$

**Lemma 3.3** ([17, Lemma 4.1]). *The stiffness matrix  $A$  and its diagonal part  $A_D$  satisfy*

$$A \leq (d+1)A_D. \quad (13)$$

**Lemma 3.4** ([15, Lemma 2.5]). *The diagonal entries of the stiffness matrix are bounded by*

$$C_{\hat{\nabla}} \sum_{K \in \omega_j} |K| \lambda_{\min} \left( (F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T} \right) \leq A_{jj} \leq C_{\hat{\nabla}} \sum_{K \in \omega_j} |K| \lambda_{\max} \left( (F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T} \right), \quad (14)$$

where  $\mathbb{D}_K = \frac{1}{|K|} \int_K \mathbb{D}(\mathbf{x}) d\mathbf{x}$  is the average of  $\mathbb{D}$  on  $K$  and  $C_{\hat{\nabla}} = \frac{d+1}{d} \left( \frac{d!}{\sqrt{d+1}} \right)^{\frac{2}{d}}$ .

The next two lemmas establish bounds for  $A_{jj}$  with a more explicit geometric interpretation than (14). For this, we denote the diameter and the minimal height of  $K$  in the metric  $\mathbb{D}_K^{-1}$  by  $h_{K, \mathbb{D}^{-1}}$  and  $a_{K, \mathbb{D}^{-1}}$ , respectively. The average element diameter is defined as

$$h_{\mathbb{D}^{-1}} = \left( \frac{|\Omega|_{\mathbb{D}^{-1}}}{N} \right)^{\frac{1}{d}}, \quad |\Omega|_{\mathbb{D}^{-1}} = \sum_{K \in \mathcal{T}_h} |K| \sqrt{\det(\mathbb{D}_K^{-1})}. \quad (15)$$

Further, let  $\hat{h}$ ,  $\hat{\rho}$ , and  $\hat{a}$  be the diameter, the in-diameter, and the minimal height of the unitary equilateral  $\hat{K}$ , respectively, i.e.,

$$\hat{h} = \sqrt{2} \left( \frac{d!}{\sqrt{d+1}} \right)^{\frac{1}{d}}, \quad \hat{\rho} = \sqrt{\frac{2}{d(d+1)}} \hat{h}, \quad \hat{a} = \sqrt{\frac{d+1}{2d}} \hat{h}.$$

**Lemma 3.5** ([13, Lemmas 4.1 and 4.2]). *It holds*

$$\frac{h_{K, \mathbb{D}^{-1}}^2}{\hat{h}^2} \leq \|(F'_K)^T \mathbb{D}_K^{-1} F'_K\|_2 \leq \frac{h_{K, \mathbb{D}^{-1}}^2}{\hat{\rho}^2}, \quad (16)$$

$$\frac{\hat{a}^2}{a_{K, \mathbb{D}^{-1}}^2} \leq \|(F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T}\|_2 \leq \frac{d^2 \hat{a}^2}{a_{K, \mathbb{D}^{-1}}^2}. \quad (17)$$

**Lemma 3.6.** *It holds*

$$\hat{\rho}^2 C_{\hat{\nabla}} h_{\mathbb{D}^{-1}}^{-2} \sum_{K \in \omega_j} |K| \left( \frac{h_{\mathbb{D}^{-1}}}{h_{K, \mathbb{D}^{-1}}} \right)^2 \leq A_{jj} \leq d^2 \hat{a}^2 C_{\hat{\nabla}} h_{\mathbb{D}^{-1}}^{-2} \sum_{K \in \omega_j} |K| \left( \frac{h_{\mathbb{D}^{-1}}}{a_{K, \mathbb{D}^{-1}}} \right)^2. \quad (18)$$

*Proof.* Combining (14) and (17) yields

$$A_{jj} \leq d^2 \hat{a}^2 C_{\hat{\nabla}} \sum_{K \in \omega_j} |K| a_{K, \mathbb{D}^{-1}}^{-2} = d^2 \hat{a}^2 C_{\hat{\nabla}} h_{\mathbb{D}^{-1}}^{-2} \sum_{K \in \omega_j} |K| \left( \frac{h_{\mathbb{D}^{-1}}}{a_{K, \mathbb{D}^{-1}}} \right)^2,$$

which gives the right inequality of (18). On the other hand, (16) gives

$$\lambda_{\min}((F'_K)^{-1} \mathbb{D}_K (F'_K)^{-T}) = \|(F'_K)^T \mathbb{D}_K^{-1} F'_K\|_2^{-1} \geq \hat{\rho}^2 h_{K, \mathbb{D}^{-1}}^{-2}.$$

Combining this with (14) leads to the left inequality of (18).  $\square$

A mesh that is uniform in the metric  $\mathbb{D}^{-1}$  (*a  $\mathbb{D}^{-1}$ -uniform mesh*) satisfies  $h_{K, \mathbb{D}^{-1}} \sim a_{K, \mathbb{D}^{-1}} \sim h_{\mathbb{D}^{-1}}$ . In such a case, (18) implies  $A_{jj} \sim |\omega_j| h_{\mathbb{D}^{-1}}^{-2} \sim |\omega_j| N^{\frac{2}{d}}$ .

**Lemma 3.7** ([17, Lemma 5.1]). *The smallest eigenvalue of the stiffness matrix is bounded by*

$$\lambda_{\min}(A) \gtrsim d_{\min} N^{-1} \Psi_E^{-1}, \quad (19)$$

where

$$\Psi_E = \begin{cases} 1, & d = 1, \\ 1 + \ln \left( \frac{|\bar{K}|}{|K_{\min}|} \right), & d = 2, \\ \left( \frac{1}{N} \sum_{K \in \mathcal{T}_h} \left( \frac{|\bar{K}|}{|K|} \right)^{\frac{d-2}{2}} \right)^{\frac{2}{d}}, & d \geq 3, \end{cases} \quad (20)$$

$|K_{\min}|$  and  $|\bar{K}| = \frac{1}{N} |\Omega|$  are the minimal and the average element volumes and  $d_{\min}$  is the global smallest eigenvalue of  $\mathbb{D}$  (cf. (2)).

The factor  $\Psi_E$  shows that the dependence of the lower bound on the mesh non-uniformity is mild although getting stronger in higher dimensions:  $\Psi_E$  is mesh-independent in 1d, contains only  $\ln(|\bar{K}|/|K_{\min}|)$  in 2d, and involves the generalized  $d/2$ -mean of the ratio  $(|\bar{K}|/|K|)^{(d-2)/d}$  over all elements for  $d \geq 3$ .

The next lemmas summarize bounds for the extremal eigenvalues of  $M$  and  $A$ .

**Lemma 3.8.** *The extremal eigenvalues of  $P^{-\frac{1}{2}} M P^{-\frac{1}{2}}$  are bounded by*

$$\lambda_{\max}(P^{-\frac{1}{2}} M P^{-\frac{1}{2}}) \leq \begin{cases} \frac{|\omega_{\max}|}{d+1}, & P = I_{N_{vi}}, \\ 1 + \frac{d}{2}, & P = M_D, M_{lump}, M_D + \Delta t A_D, \end{cases} \quad (21)$$

and

$$\lambda_{\min}(P^{-\frac{1}{2}} M P^{-\frac{1}{2}}) \geq \begin{cases} \frac{|\omega_{\min}|}{(d+1)(d+2)}, & P = I_{N_{vi}}, \\ 1/2, & P = M_D, \\ 1/(d+2), & P = M_{lump}, \\ C \left( 1 + \Delta t h_{\mathbb{D}^{-1}}^{-2} \max_j \sum_{K \in \omega_j} \frac{|K|}{|\omega_j|} \left( \frac{h_{\mathbb{D}^{-1}}}{a_{K, \mathbb{D}^{-1}}} \right)^2 \right)^{-1}, & P = M_D + \Delta t A_D, \end{cases} \quad (22)$$

where  $|\omega_{\max}|$  and  $|\omega_{\min}|$  are the maximal and minimal patch volumes, respectively.

*Proof.* The inequalities (21) and (22) for  $P = I_{N_{vi}}$  (no preconditioning),  $P = M_D$ , and  $M_{lump}$  follow from Lemmas 3.1 and 3.2. For  $P = M_D + \Delta t A_D$ , (11) and (18) give

$$P_{jj} \lesssim |\omega_j| + \Delta t h_{\mathbb{D}^{-1}}^{-2} \sum_{K \in \omega_j} |K| \left( \frac{h_{\mathbb{D}^{-1}}}{a_{K, \mathbb{D}^{-1}}} \right)^2. \quad (23)$$

Then, (22) follows from

$$\mathbf{v}^T P^{-\frac{1}{2}} M P^{-\frac{1}{2}} \mathbf{v} \gtrsim \sum_j \frac{v_j^2 |\omega_j|}{P_{jj}} \gtrsim \sum_j v_j^2 \left( 1 + \Delta t h_{\mathbb{D}^{-1}}^{-2} \sum_{K \in \omega_j} \frac{|K|}{|\omega_j|} \left( \frac{h_{\mathbb{D}^{-1}}}{a_{K, \mathbb{D}^{-1}}} \right)^2 \right)^{-1}. \quad \square$$

**Lemma 3.9.** *The extremal eigenvalues of  $P^{-\frac{1}{2}} A P^{-\frac{1}{2}}$  are bounded by*

$$\lambda_{\max}(P^{-\frac{1}{2}} A P^{-\frac{1}{2}}) \leq \begin{cases} Ch_{\mathbb{D}^{-1}}^{-2} \max_j \sum_{K \in \omega_j} |K| \left( \frac{h_{\mathbb{D}^{-1}}}{a_{K, \mathbb{D}^{-1}}} \right)^2, & P = I_{N_{vi}}, \\ Ch_{\mathbb{D}^{-1}}^{-2} \max_j \sum_{K \in \omega_j} \frac{|K|}{|\omega_j|} \left( \frac{h_{\mathbb{D}^{-1}}}{a_{K, \mathbb{D}^{-1}}} \right)^2, & P = M, M_D, M_{lump}, \\ (d+1)\Delta t^{-1}, & P = M_D + \Delta t A_D, \end{cases} \quad (24)$$

and

$$\lambda_{\min}(P^{-\frac{1}{2}} A P^{-\frac{1}{2}}) \gtrsim \begin{cases} d_{\min} N^{-1} \Psi_E^{-1}, & P = I_{N_{vi}}, \\ \lambda_{\mathbb{D}}, & P = M, M_D, M_{lump}, \\ d_{\min} \Psi_{\mathbb{D}}^{-1}, & P = M_D + \Delta t A_D, \end{cases} \quad (25)$$

where  $\lambda_{\mathbb{D}}$  is the minimal eigenvalue of the operator  $-\nabla \cdot (\mathbb{D} \nabla)$  and

$$\Psi_{\mathbb{D}} = \begin{cases} 1 + \Delta t h_{\mathbb{D}^{-1}}^{-2} \sum_K |K| \left( \frac{h_{\mathbb{D}^{-1}}}{a_{K, \mathbb{D}^{-1}}} \right)^2, & d = 1, \\ (1 + |\ln \psi_{\mathbb{D}}|) \left( 1 + \Delta t h_{\mathbb{D}^{-1}}^{-2} \sum_K |K| \left( \frac{h_{\mathbb{D}^{-1}}}{a_{K, \mathbb{D}^{-1}}} \right)^2 \right), & d = 2, \\ 1 + \Delta t h_{\mathbb{D}^{-1}}^{-2} \left( \sum_K |K| \left( \frac{h_{\mathbb{D}^{-1}}}{a_{K, \mathbb{D}^{-1}}} \right)^d \right)^{\frac{2}{d}}, & d \geq 3, \end{cases} \quad (26)$$

$$\psi_{\mathbb{D}} = \frac{1 + \Delta t h_{\mathbb{D}^{-1}}^{-2} \max_K \left( \frac{h_{\mathbb{D}^{-1}}}{a_{K, \mathbb{D}^{-1}}} \right)^2}{1 + \Delta t h_{\mathbb{D}^{-1}}^{-2} \sum_K \frac{|K|}{|\Omega|} \left( \frac{h_{\mathbb{D}^{-1}}}{a_{K, \mathbb{D}^{-1}}} \right)^2}. \quad (27)$$

*Proof.* For  $P = I_{N_{vi}}$ , (24) follows from Lemmas 3.3 and 3.6. For  $P = M$ , (11) and (13) give

$$\lambda_{\max}(P^{-\frac{1}{2}} A P^{-\frac{1}{2}}) = \max_{\mathbf{v} \neq 0} \frac{\mathbf{v}^T A \mathbf{v}}{\mathbf{v}^T P \mathbf{v}} \leq 2(d+1) \max_{\mathbf{v} \neq 0} \frac{\mathbf{v}^T A_D \mathbf{v}}{\mathbf{v}^T M_D \mathbf{v}} = 2(d+1) \max_j \frac{A_{jj}}{M_{jj}},$$

which, together with (11) and (18), implies (24). For  $P = M_D$  and  $P = M_{lump}$ , (24) follows from Lemmas 3.1 and 3.2. For  $P = M_D + \Delta t A_D$ , (24) follows from (13).

The inequality (25) for  $P = I_{N_{vi}}$  follows from Lemma 3.7. For  $P = M$  and, similarly for  $P = M_D$  and  $P = M_{lump}$ , (25) follows from the basic property of the conformal FE approximation of elliptic eigenvalue problems (e.g., [6, Theorem 1]), since the eigenvalue problem for  $P^{-\frac{1}{2}} A P^{-\frac{1}{2}}$  is a FE approximation to the eigenvalue problem for the operator  $-\nabla \cdot (\mathbb{D} \nabla)$ .

For  $P = M_D + \Delta t A_D$ , using the strategy in the proof of [17, Lemma 5.1], we have

$$\mathbf{v}^T P^{-\frac{1}{2}} A P^{-\frac{1}{2}} \mathbf{v} \gtrsim d_{\min} \times \begin{cases} \left( \sum_j v_j^2 \right) \cdot \left( \sum_j P_{jj} \right)^{-1}, & d = 1, \\ \frac{1}{q} \left( \sum_K s_K^{\frac{q}{q-2}} \right)^{-\frac{q-2}{q}} \sum_j v_j^2 P_{jj}^{-1} \sum_{K \in \omega_j} s_K |K|^{\frac{2}{q}}, & d = 2, \\ \left( \sum_K s_K^{\frac{d}{2}} \right)^{-\frac{2}{d}} \sum_j v_j^2 P_{jj}^{-1} \sum_{K \in \omega_j} s_K |K|^{\frac{d-2}{d}}, & d \geq 3, \end{cases} \quad (28)$$

with a parameter  $q > 2$  and some not-all-zero nonnegative numbers  $s_K, K \in \mathcal{T}_h$  (to be determined later). Here,  $\sum_K$  denotes the summation over all elements in  $\mathcal{T}_h$ .

For notational simplicity, we denote  $r_K = a_{K, \mathbb{D}-1}^{-2}$  and rewrite (23) as

$$P_{jj} \lesssim \sum_{K \in \omega_j} |K| (1 + \Delta t r_K). \quad (29)$$

For  $d = 1$ , (29) implies

$$\mathbf{v}^T P^{-\frac{1}{2}} A P^{-\frac{1}{2}} \mathbf{v} \gtrsim \frac{d_{\min} \sum_j v_j^2}{\sum_K |K| (1 + \Delta t r_K)}$$

and therefore

$$\lambda_{\min}(P^{-\frac{1}{2}} A P^{-\frac{1}{2}}) \gtrsim \frac{d_{\min}}{\sum_K |K| (1 + \Delta t r_K)}. \quad (30)$$

For  $d = 2$ , we choose  $s_K = |K|^{1-\frac{2}{q}} (1 + \Delta t r_K)$ . Then, (29) implies

$$P_{jj}^{-1} \sum_{K \in \omega_j} s_K |K|^{\frac{2}{q}} \geq C$$

and using (28) we obtain

$$\mathbf{v}^T P^{-\frac{1}{2}} A P^{-\frac{1}{2}} \mathbf{v} \gtrsim \frac{d_{\min} \sum_k v_k^2}{q \left( \sum_K |K| (1 + \Delta t r_K)^{\frac{q}{q-2}} \right)^{\frac{q-2}{q}}}.$$

The denominator tends to infinity as  $q \rightarrow \infty$  and to  $2 \max_K (1 + \Delta t r_K)$  as  $q \rightarrow 2^+$ . To find an optimal choice for  $q$ , we first estimate the denominator as

$$\begin{aligned} q \left( \sum_K |K| (1 + \Delta t r_K)^{\frac{q}{q-2}} \right)^{\frac{q-2}{q}} &= q \left( \sum_K |K| (1 + \Delta t r_K) \cdot (1 + \Delta t r_K)^{\frac{2}{q-2}} \right)^{\frac{q-2}{q}} \\ &\leq q \left( \left( \sum_K |K| (1 + \Delta t r_K) \right) \cdot \left( 1 + \Delta t \max_K r_K \right)^{\frac{2}{q-2}} \right)^{\frac{q-2}{q}} \\ &= q \psi_{\mathbb{D}}^{\frac{2}{q}} \sum_K |K| (1 + \Delta t r_K), \end{aligned}$$



where

$$\psi_{\mathbb{D}} = \frac{1 + \Delta t \max_K r_K}{\sum_K |K| (1 + \Delta t r_K)}. \quad (31)$$

If  $\psi_{\mathbb{D}} > e$ , we choose  $q = 2 \ln \psi_{\mathbb{D}}$  and obtain  $q\psi_{\mathbb{D}}^{\frac{2}{q}} = 2e \ln \psi_{\mathbb{D}}$ . If  $\psi_{\mathbb{D}} \leq e$ , we use  $q \rightarrow 2^+$  and obtain  $q\psi_{\mathbb{D}}^{\frac{2}{q}} \leq 2e$ . Combining these two cases yields

$$q \left( \sum_K |K| (1 + \Delta t r_K)^{\frac{q}{q-2}} \right)^{\frac{q-2}{q}} \leq 2e (1 + |\ln \psi_{\mathbb{D}}|) \left( \sum_K |K| (1 + \Delta t r_K) \right)$$

and therefore

$$\lambda_{\min}(P^{-\frac{1}{2}} A P^{-\frac{1}{2}}) \gtrsim \frac{d_{\min}}{(1 + |\ln \psi_{\mathbb{D}}|) \left( \sum_K |K| (1 + \Delta t r_K) \right)}. \quad (32)$$

For  $d \geq 3$ , we choose  $s_K = |K|^{\frac{2}{d}} (1 + \Delta t r_K)$  such that  $P_{jj}^{-1} \sum_{K \in \omega_j} r_K |K|^{\frac{2}{d}} \geq C$ . From (28), we have

$$\mathbf{v}^T P^{-\frac{1}{2}} A P^{-\frac{1}{2}} \mathbf{u} \gtrsim \frac{d_{\min} \sum_j v_j^2}{\left( \sum_K |K| (1 + \Delta t r_K)^{\frac{d}{2}} \right)^{\frac{2}{d}}},$$

which gives

$$\lambda_{\min}(P^{-\frac{1}{2}} A P^{-\frac{1}{2}}) \gtrsim \frac{d_{\min}}{\left( \sum_K |K| (1 + \Delta t r_K)^{\frac{d}{2}} \right)^{\frac{2}{d}}}. \quad (33)$$

Further, (31) can be rewritten into (27) and

$$\begin{aligned} \sum_K |K| (1 + \Delta t r_K) &= |\Omega| + \Delta t h_{\mathbb{D}^{-1}}^{-2} \sum_K |K| \left( \frac{h_{\mathbb{D}^{-1}}}{a_{K, \mathbb{D}^{-1}}} \right)^2, \\ \left( \sum_K |K| (1 + \Delta t r_K)^{\frac{d}{2}} \right)^{\frac{2}{d}} &\lesssim |\Omega|^{\frac{2}{d}} + \Delta t h_{\mathbb{D}^{-1}}^{-2} \left( \sum_K |K| \left( \frac{h_{\mathbb{D}^{-1}}}{a_{K, \mathbb{D}^{-1}}} \right)^d \right)^{\frac{2}{d}}. \end{aligned}$$

Combining these with (30) to (33) gives (25) with  $P = M_D + \Delta t A_D$ .  $\square$

## 4 Conditioning of the implicit Euler integration

In the following we estimate the conditioning of  $M + \Delta t A$  (in the  $l_2$ -norm), which is related to the efficient iterative solution of (9).

### 4.1 Condition number of $M + \Delta t A$

**Theorem 4.1.** *The condition number of  $M + \Delta t A$  is bounded by*

$$\kappa(M + \Delta t A) \lesssim \frac{\max_j \frac{|\omega_j|}{|\omega_{\min}|} \left( 1 + \Delta t h_{\mathbb{D}^{-1}}^{-2} \sum_{K \in \omega_j} \frac{|K|}{|\omega_j|} \left( \frac{h_{\mathbb{D}^{-1}}}{a_{K, \mathbb{D}^{-1}}} \right)^2 \right)}{1 + \Delta t d_{\min} N^{-1} |\omega_{\min}|^{-1} \Psi_E^{-1}}. \quad (34)$$

*Proof.* (11) and (13) yield

$$M + \Delta t A \leq \frac{d+2}{2} M_D + \Delta t (d+1) A_D \leq (d+1)(M_D + \Delta t A_D).$$

From this and (11) and (18) we obtain

$$\lambda_{\max}(M + \Delta t A) \lesssim \max_j \left( |\omega_j| + \Delta t h_{\mathbb{D}^{-1}}^{-2} \sum_{K \in \omega_j} |K| \left( \frac{h_{\mathbb{D}^{-1}}}{a_{K, \mathbb{D}^{-1}}} \right)^2 \right).$$

On the other hand, using  $\lambda_{\min}(M + \Delta t A) \geq \lambda_{\min}(M) + \Delta t \lambda_{\min}(A)$  and Lemmas 3.8 and 3.9 we have

$$\lambda_{\min}(M + \Delta t A) \gtrsim |\omega_{\min}| + \Delta t d_{\min} N^{-1} \Psi_E^{-1}.$$

Combining the above results yields (34).  $\square$

There are three factors influencing bound (34). The first factor is the number of the mesh elements  $N$ . The second factor is

$$\sum_{K \in \omega_j} \frac{|K|}{|\omega_j|} \left( \frac{h_{\mathbb{D}^{-1}}}{a_{K, \mathbb{D}^{-1}}} \right)^2, \quad (35)$$

which reflects the mesh nonuniformity in the metric  $\mathbb{D}^{-1}$  and is a constant for a  $\mathbb{D}^{-1}$ -uniform mesh, for which  $a_{K, \mathbb{D}^{-1}} \sim h_{\mathbb{D}^{-1}}$  for all  $K$ . The third factor is the effect of the mesh nonuniformity (in the Euclidean metric) reflected by  $\Psi_E$ ,  $|\omega_j|/|\omega_{\min}|$ , and  $N|\omega_{\min}|$ , which all become constants if the mesh is uniform.

The time step size  $\Delta t$  plays the role of a homotopy parameter between the mass and the stiffness matrices:

$$\kappa(M + \Delta t A) \xrightarrow{\Delta t \rightarrow 0} \kappa(M) \lesssim \frac{|\omega_{\max}|}{|\omega_{\min}|}. \quad (36)$$

Thus,  $\kappa(M)$  depends on the nonuniformity of the element patch volumes. On the other hand,

$$\kappa(M + \Delta t A) \xrightarrow{\Delta t \rightarrow \infty} \kappa(A) \lesssim d_{\min}^{-1} h_{\mathbb{D}^{-1}}^{-2} \Psi_E \max_j \left( N \sum_{K \in \omega_j} |K| \left( \frac{h_{\mathbb{D}^{-1}}}{a_{K, \mathbb{D}^{-1}}} \right)^2 \right), \quad (37)$$

which has been obtained previously [17, Theorem 5.2]. It depends on three factors as well:  $N$  (through  $h_{\mathbb{D}^{-1}}$ ), the mesh nonuniformity in the Euclidean metric (through  $\Psi_E$ ), and the mesh nonuniformity with respect to  $\mathbb{D}^{-1}$  (through the term in the maximum).

## 4.2 Condition number of $P^{-\frac{1}{2}}(M + \Delta t A)P^{-\frac{1}{2}}$

**Theorem 4.2.** *The condition number of  $P^{-\frac{1}{2}}(M + \Delta t A)P^{-\frac{1}{2}}$  with  $P = M$ ,  $M_D$ , or  $M_{lump}$  is bounded by*

$$\kappa(P^{-\frac{1}{2}}(M + \Delta t A)P^{-\frac{1}{2}}) \lesssim \frac{1 + \Delta t h_{\mathbb{D}^{-1}}^{-2} \max_j \sum_{K \in \omega_j} \frac{|K|}{|\omega_j|} \left( \frac{h_{\mathbb{D}^{-1}}}{a_{K, \mathbb{D}^{-1}}} \right)^2}{1 + \Delta t \lambda_{\mathbb{D}}}, \quad (38)$$

where  $\lambda_{\mathbb{D}}$  is the minimal eigenvalue of  $-\nabla \cdot (\mathbb{D} \nabla)$ .

*Proof.* Bound (38) is obtained by combining Lemmas 3.8 and 3.9 with the estimates

$$\begin{aligned}\lambda_{\max} \left( P^{-\frac{1}{2}}(M + \Delta t A)P^{-\frac{1}{2}} \right) &\leq \lambda_{\max} \left( P^{-\frac{1}{2}}MP^{-\frac{1}{2}} \right) + \Delta t \lambda_{\max} \left( P^{-\frac{1}{2}}AP^{-\frac{1}{2}} \right), \\ \lambda_{\min} \left( P^{-\frac{1}{2}}(M + \Delta t A)P^{-\frac{1}{2}} \right) &\geq \lambda_{\min} \left( P^{-\frac{1}{2}}MP^{-\frac{1}{2}} \right) + \Delta t \lambda_{\min} \left( P^{-\frac{1}{2}}AP^{-\frac{1}{2}} \right).\end{aligned}\quad \square$$

The bounds on  $\kappa(P^{-\frac{1}{2}}(M + \Delta t A)P^{-\frac{1}{2}})$  for  $P = M$ ,  $M_D$ , and  $M_{lump}$  are similar, while the last two choices lead to a simple diagonal scaling, which is easier to implement.

All three choices reduce the effects of the mesh nonuniformity: in comparison to (34), bound (38) does not depend on  $|\omega_{\min}|$  or  $\Psi_E$  directly and contains only the  $\mathbb{D}^{-1}$ -nonuniformity factor (35). This is intuitive, since the eigenvalues of  $M^{-1}A$  approximate those of the underlying continuous operator, which are mesh-independent. However,  $\kappa(M^{-1}A)$  is not necessarily smaller than  $\kappa(A)$  and the overall effect depends on the magnitude of  $\Delta t$ : if  $\Delta t$  is large,  $\kappa(M^{-1}(M + \Delta t A))$  might not be better than  $\kappa(M + \Delta t A)$ . On the other hand, if  $\Delta t$  is small, we can expect that  $\kappa(M^{-1}(M + \Delta t A)) < \kappa(M + \Delta t A)$ . The numerical experiments in Sect. 6 support this argument (see Example 6.2 and Fig. 3b).

**Theorem 4.3.** *The condition number of  $P^{-\frac{1}{2}}(M + \Delta t A)P^{-\frac{1}{2}}$  with the Jacobi preconditioner  $P = M_D + \Delta t A_D$  is bounded by*

$$\kappa(P^{-\frac{1}{2}}(M + \Delta t A)P^{-\frac{1}{2}}) \lesssim \left( \frac{1}{1 + \Delta t h_{\mathbb{D}^{-1}}^{-2} \max_j \sum_{K \in \omega_j} \frac{|K|}{|\omega_j|} \left( \frac{h_{\mathbb{D}^{-1}}}{a_{K, \mathbb{D}^{-1}}} \right)^2} + \frac{d_{\min}}{\Psi_{\mathbb{D}}} \right)^{-1} \quad (39)$$

where  $\Psi_{\mathbb{D}}$  is given in (26).

*Proof.* The proof is similar to that of Theorem 4.2.  $\square$

Bound (39) is comparable to (38) although the former is smaller than the latter in general, especially for large  $\Delta t$ , since the factor  $\Psi_{\mathbb{D}}$  in (39) involves averaging over all elements, whereas (38) involves the maximum over patch averages.

In (39),  $\Delta t$  is a homotopy parameter between the mass and stiffness matrices:

$$\kappa(P^{-\frac{1}{2}}(M + \Delta t A)P^{-\frac{1}{2}}) \xrightarrow{\Delta t \rightarrow 0} \kappa(M_D^{-\frac{1}{2}}MM_D^{-\frac{1}{2}}) \leq C.$$

On the other hand,

$$\begin{aligned}\kappa(P^{-\frac{1}{2}}(M + \Delta t A)P^{-\frac{1}{2}}) &\xrightarrow{\Delta t \rightarrow \infty} \kappa(A_D^{-\frac{1}{2}}AA_D^{-\frac{1}{2}}) \\ &\lesssim \frac{N_d^{\frac{2}{d}}}{d_{\min}|\Omega|_{\mathbb{D}^{-1}}^{\frac{2}{d}}} \times \begin{cases} \sum_K |K| \left( \frac{h_{\mathbb{D}^{-1}}}{a_{K, \mathbb{D}^{-1}}} \right)^2, & d = 1, \\ (1 + |\ln \psi_{\mathbb{D}}|) \sum_K |K| \left( \frac{h_{\mathbb{D}^{-1}}}{a_{K, \mathbb{D}^{-1}}} \right)^2, & d = 2, \\ \left( \sum_K |K| \left( \frac{h_{\mathbb{D}^{-1}}}{a_{K, \mathbb{D}^{-1}}} \right)^d \right)^{\frac{2}{d}}, & d \geq 3, \end{cases}\end{aligned}$$

where  $\psi_{\mathbb{D}}$  becomes

$$\psi_{\mathbb{D}} = \frac{\max_K \left( \frac{h_{\mathbb{D}^{-1}}}{a_{K, \mathbb{D}^{-1}}} \right)^2}{\sum_K \frac{|K|}{|\Omega|} \left( \frac{h_{\mathbb{D}^{-1}}}{a_{K, \mathbb{D}^{-1}}} \right)^2}.$$

This bound is equivalent to the bound obtained in [17, Theorem 5.2].

*Remark 4.1.* In comparison to the bound (34) for  $M + \Delta t A$ , bounds (38) and (39) contain only two mesh-dependent factors:  $N$  and the mesh nonuniformity in the metric  $\mathbb{D}^{-1}$  (through the terms involving the ratio  $h_{\mathbb{D}^{-1}}/a_{K,\mathbb{D}^{-1}}$ ). This shows that the effects of the mesh nonuniformity (in the Euclidean metric) on the condition number is effectively eliminated by the preconditioning.

*Remark 4.2.* For a  $\mathbb{D}^{-1}$ -uniform mesh,  $a_{K,\mathbb{D}^{-1}} \sim h_{\mathbb{D}^{-1}}$  and, hence, all terms involving the ratio  $h_{\mathbb{D}^{-1}}/a_{K,\mathbb{D}^{-1}}$  will become a constant. Thus,

$$\kappa(P^{-\frac{1}{2}}(M + \Delta t A)P^{-\frac{1}{2}}) = \mathcal{O}(1 + \Delta t N^{\frac{2}{d}}),$$

which shows more clearly the role of  $\Delta t$ :

$$\kappa(P^{-\frac{1}{2}}(M + \Delta t A)P^{-\frac{1}{2}}) = \begin{cases} \mathcal{O}(N^{\frac{2}{d}}), & \Delta t = \mathcal{O}(1), \\ \mathcal{O}(N^{\frac{1}{d}}), & \Delta t = \mathcal{O}(h_{\mathbb{D}^{-1}}) = \mathcal{O}(N^{-\frac{1}{d}}), \\ \mathcal{O}(1), & \Delta t = \mathcal{O}(N^{-\frac{2}{d}}) \text{ or } \Delta t \rightarrow 0. \end{cases}$$

## 5 Conditioning of general implicit RK integration

For a general implicit RK method, the matrix (8) is not necessarily normal and we have to work with singular values for its condition number. We first recall some properties of singular values and condition numbers of matrices. In the following,  $\|\cdot\|$  always denotes the  $l_2$ -norm.

For an arbitrary square matrix  $U$ , its maximal and minimal singular values are

$$\sigma_{\max}(U) = \max_{\mathbf{v} \neq 0} \frac{\|U\mathbf{v}\|}{\|\mathbf{v}\|} \quad \text{and} \quad \sigma_{\min}(U) = \min_{\mathbf{v} \neq 0} \frac{\|U\mathbf{v}\|}{\|\mathbf{v}\|}$$

and the condition number of  $U$  is defined as

$$\kappa(U) = \frac{\sigma_{\max}(U)}{\sigma_{\min}(U)}.$$

If  $U$  is normal, it becomes

$$\kappa(U) = \frac{|\lambda_{\max}(U)|}{|\lambda_{\min}(U)|}.$$

Further, for any square matrices  $U$  and  $V$  we have

$$\sigma_{\max}(U + V) \leq \sigma_{\max}(U) + \sigma_{\max}(V). \quad (40)$$

The analogue for the minimal singular value does not hold in general. However, the following lemma provides a sufficient condition.

**Lemma 5.1.** *If the eigenvalues of  $U^T V$  have nonnegative real parts, then*

$$\sigma_{\min}^2(U + V) \geq \sigma_{\min}^2(U) + \sigma_{\min}^2(V). \quad (41)$$

*Proof.*

$$\begin{aligned} \sigma_{\min}^2(U + V) &= \min_{\mathbf{v} \neq 0} \frac{\|(U + V)\mathbf{v}\|_2^2}{\|\mathbf{v}\|^2} = \min_{\mathbf{v} \neq 0} \frac{\mathbf{v}^T (U^T + V^T)(U + V)\mathbf{v}}{\|\mathbf{v}\|^2} \\ &\geq \min_{\mathbf{v} \neq 0} \frac{\mathbf{v}^T U^T U \mathbf{v}}{\|\mathbf{v}\|^2} + \min_{\mathbf{v} \neq 0} \frac{\mathbf{v}^T V^T V \mathbf{v}}{\|\mathbf{v}\|^2} + \min_{\mathbf{v} \neq 0} \frac{\mathbf{v}^T (U^T V + V^T U)\mathbf{v}}{\|\mathbf{v}\|^2} \\ &= \sigma_{\min}^2(U) + \sigma_{\min}^2(V) + \min_{\mathbf{v} \neq 0} \frac{\mathbf{v}^T (U^T V + V^T U)\mathbf{v}}{\|\mathbf{v}\|^2}. \end{aligned}$$

Since the eigenvalues of  $U^T V$  have nonnegative real parts, the third term on the right-hand side of the above equation is nonnegative and (41) follows.  $\square$

The next lemma provides the eigenvalues and singular values of the Kronecker product of any square matrices  $U$  and  $V$ .

**Lemma 5.2** ([18, Theorems 13.10 and 13.12]). *For any square matrices  $U$  and  $V$ , the eigenvalues and singular values of  $U \otimes V$  are  $\lambda_j(U)\lambda_k(V)$  and  $\sigma_j(U)\sigma_k(V)$ ,  $j, k = 1, 2, \dots$*

Lemma 5.2 in particular implies

$$\sigma_{\max}(U \otimes V) = \sigma_{\max}(U)\sigma_{\max}(V) \quad \text{and} \quad \sigma_{\min}(U \otimes V) = \sigma_{\min}(U)\sigma_{\min}(V). \quad (42)$$

## 5.1 General implicit RK methods: simultaneous solution

Although a successive strategy for implicit RK methods is more common (e.g., [9, p. 131]), the simultaneous solution can be appealing if an iterative method is used. It is also important to have a theoretical understanding of the overall system conditioning.

**Theorem 5.1.** *Assume that the eigenvalues of the coefficient matrix  $\Gamma$  for a given implicit RK method have nonnegative real parts and let  $P$  be a symmetric and positive definite preconditioner. Then,*

$$\begin{aligned} \kappa \left( (I_s \otimes P^{-\frac{1}{2}})(I_s \otimes M + \Delta t \Gamma \otimes A)(I_s \otimes P^{-\frac{1}{2}}) \right) \\ \leq \frac{\lambda_{\max}(P^{-\frac{1}{2}} M P^{-\frac{1}{2}}) + \Delta t \sigma_{\max}(\Gamma) \lambda_{\max}(P^{-\frac{1}{2}} A P^{-\frac{1}{2}})}{\sqrt{\lambda_{\min}^2(P^{-\frac{1}{2}} M P^{-\frac{1}{2}}) + \Delta t^2 \sigma_{\min}^2(\Gamma) \lambda_{\min}^2(P^{-\frac{1}{2}} A P^{-\frac{1}{2}})}}. \end{aligned} \quad (43)$$

*Proof.* First, notice that

$$(I_s \otimes P^{-\frac{1}{2}})(I_s \otimes M + \Delta t \Gamma \otimes A)(I_s \otimes P^{-\frac{1}{2}}) = I_s \otimes (P^{-\frac{1}{2}} M P^{-\frac{1}{2}}) + \Delta t \Gamma \otimes (P^{-\frac{1}{2}} A P^{-\frac{1}{2}}).$$

Then, from (40) and (42) we have

$$\begin{aligned} \sigma_{\max} \left( (I_s \otimes P^{-\frac{1}{2}})(I_s \otimes M + \Delta t \Gamma \otimes A)(I_s \otimes P^{-\frac{1}{2}}) \right) \\ \leq \sigma_{\max} \left( I_s \otimes (P^{-\frac{1}{2}} M P^{-\frac{1}{2}}) \right) + \Delta t \sigma_{\max} \left( \Gamma \otimes (P^{-\frac{1}{2}} A P^{-\frac{1}{2}}) \right) \\ = \sigma_{\max} \left( P^{-\frac{1}{2}} M P^{-\frac{1}{2}} \right) + \Delta t \sigma_{\max}(\Gamma) \sigma_{\max} \left( P^{-\frac{1}{2}} A P^{-\frac{1}{2}} \right) \\ = \lambda_{\max} \left( P^{-\frac{1}{2}} M P^{-\frac{1}{2}} \right) + \Delta t \sigma_{\max}(\Gamma) \lambda_{\max} \left( P^{-\frac{1}{2}} A P^{-\frac{1}{2}} \right). \end{aligned} \quad (44)$$

Moreover,

$$\begin{aligned} \left( I_s \otimes (P^{-\frac{1}{2}} M P^{-\frac{1}{2}}) \right)^T \left( \Gamma \otimes (P^{-\frac{1}{2}} A P^{-\frac{1}{2}}) \right) \\ = \left( I_s \otimes (P^{-\frac{1}{2}} M P^{-\frac{1}{2}}) \right) \left( \Gamma \otimes (P^{-\frac{1}{2}} A P^{-\frac{1}{2}}) \right) = \Gamma \otimes (P^{-\frac{1}{2}} M P^{-1} A P^{-\frac{1}{2}}). \end{aligned} \quad (45)$$

$P^{-\frac{1}{2}} M P^{-1} A P^{-\frac{1}{2}}$  has real positive eigenvalues and the eigenvalues of  $\Gamma$  have nonnegative real parts by assumption. Hence, Lemma 5.2 implies that the eigenvalues of the matrix (45) have nonnegative real parts as well. Using Lemma 5.1 we get

$$\begin{aligned} \sigma_{\min}^2 \left( (I_s \otimes P^{-\frac{1}{2}})(I_s \otimes M + \Delta t \Gamma \otimes A)(I_s \otimes P^{-\frac{1}{2}}) \right) \\ \geq \lambda_{\min}^2 \left( P^{-\frac{1}{2}} M P^{-\frac{1}{2}} \right) + \Delta t^2 \sigma_{\min}^2(\Gamma) \lambda_{\min}^2 \left( P^{-\frac{1}{2}} A P^{-\frac{1}{2}} \right). \end{aligned}$$

Combining this with (44) results in (43).  $\square$

Table 1: Eigenvalues and singular values of the implicit RK coefficient matrix  $\Gamma$ .

Method	Order	Eigenvalues	$\sigma_{\max}$	$\sigma_{\min}$
Gauss	4	$0.25 \pm 0.1443i$	0.6319	0.1319
Gauss	6	$0.2153, 0.1423 \pm 0.1358i$	0.6629	0.0635
Radau IA	3	$0.3333 \pm 0.2357i$	0.5	0.3333
Radau IIA	5	$0.2749, 0.1626 \pm 0.1849i$	0.8023	0.0923
Lobatto IIIA	4	$0, 0.25 \pm 0.1443i$	0.7947	0

*Remark 5.1.* Since  $\sigma_{\max}(\Gamma)$  and  $\sigma_{\min}(\Gamma)$  are typically  $\mathcal{O}(1)$  for implicit Runge-Kutta methods (see Table 1 for examples), Theorem 5.1 implies that the conditioning of the system (8) resulting from implicit RK integration behaves just like that of  $M + \Delta t A$  for the implicit Euler method. Hence, the bounds also involve three mesh-dependent factors and one of them can be eliminated effectively by diagonal preconditioning (see Remark 4.1).

In case of  $P = I_{N_{vi}}$ ,  $M$ ,  $M_D$ ,  $M_{lump}$ , and  $M_D + \Delta t A_D$ , the estimates for  $\lambda_{\min}(P^{-\frac{1}{2}}MP^{-\frac{1}{2}})$ ,  $\lambda_{\max}(P^{-\frac{1}{2}}MP^{-\frac{1}{2}})$ ,  $\lambda_{\min}(P^{-\frac{1}{2}}AP^{-\frac{1}{2}})$ , and  $\lambda_{\max}(P^{-\frac{1}{2}}AP^{-\frac{1}{2}})$  are given by Lemmas 3.8 and 3.9.

## 5.2 Diagonally implicit RK (DIRK) methods

The coefficient matrix  $\Gamma = (\gamma_{kj})_{k,j=1}^s$  of a DIRK method is a lower triangular matrix and the system (6) is solved by successively solving  $s$  linear systems with  $M + \Delta t \gamma_{jj}A$ ,  $j = 1, \dots, s$ . The conditioning of  $M + \Delta t \gamma_{jj}A$  is therefore similar to that of  $M + \Delta t A$  and, hence, the analysis in Sect. 4 also applies to DIRK methods.

## 5.3 General implicit RK methods: successive solution

A successive solution procedure transforms the large system (6) into a number of smaller systems, which are then solved successively. We adopt the approach of Butcher [4] and Bickart [3] and carry out the transformation using the Jordan normal form of the RK matrix. To keep our analysis applicable to methods with a singular  $\Gamma$ , we use the Jordan normal form of  $\Gamma$  instead of  $\Gamma^{-1}$ , which is the conventional choice (e.g., see [9, p. 131]).

Let the Jordan normal form of  $\Gamma$  be

$$\Gamma = T \begin{bmatrix} J_1 & & \\ & \ddots & \\ & & J_p \end{bmatrix} T^{-1}, \quad (46)$$

where  $T$  is a real invertible matrix and  $J_j$ ,  $j = 1, \dots, p$ , are the Jordan blocks, which either have the form

$$J_j = \begin{bmatrix} \mu_j & 1 & & \\ & \mu_j & \ddots & \\ & & \ddots & 1 \\ & & & \mu_j \end{bmatrix}, \quad \mu_j \in \mathbb{R}, \quad (47)$$

or

$$J_j = \begin{bmatrix} C_j & I_2 & & \\ & C_j & \ddots & \\ & & \ddots & I_2 \\ & & & C_j \end{bmatrix} \quad \text{with} \quad C_j = \begin{bmatrix} \alpha_j & \beta_j \\ -\beta_j & \alpha_j \end{bmatrix}, \quad \alpha_j, \beta_j \in \mathbb{R}. \quad (48)$$

Recall that the eigenvalues of  $\Gamma$  are assumed to have nonnegative real parts, i.e.,  $\mu_j, \alpha_j \geq 0$ .

**Theorem 5.2.** *Assume that the eigenvalues of the coefficient matrix  $\Gamma$  for a given implicit RK method have nonnegative real parts and the system (6) is solved by using the Jordan normal form of  $\Gamma$  to transform it into smaller systems. Then, the conditioning of implicit RK integration of (4) is determined by the conditioning of*

$$M + \mu_j \Delta t A \quad \text{and} \quad I_2 \otimes M + \Delta t C_j \otimes A, \quad (49)$$

where  $\mu_j$  and  $\alpha_j \pm i\beta_j$  are the real and complex eigenvalues of  $\Gamma$ , respectively.

Moreover, for any symmetric and positive definite preconditioner  $P$ ,

$$\begin{aligned} \kappa \left( (I_2 \otimes P^{-\frac{1}{2}})(I_2 \otimes M + \Delta t C_j \otimes A)(I_2 \otimes P^{-\frac{1}{2}}) \right) \\ \leq \frac{\lambda_{\max}(P^{-\frac{1}{2}} M P^{-\frac{1}{2}}) + \Delta t \sqrt{\alpha_j^2 + \beta_j^2} \lambda_{\max}(P^{-\frac{1}{2}} A P^{-\frac{1}{2}})}{\sqrt{\lambda_{\min}^2(P^{-\frac{1}{2}} M P^{-\frac{1}{2}}) + \Delta t^2(\alpha_j^2 + \beta_j^2) \lambda_{\min}^2(P^{-\frac{1}{2}} A P^{-\frac{1}{2}})}}. \end{aligned} \quad (50)$$

*Proof.* Using (46), we can rewrite (8) as

$$\begin{aligned} \mathcal{A} &= (T \otimes I_{N_{vi}}) \left( I_s \otimes M + \Delta t \begin{bmatrix} J_1 & & \\ & \ddots & \\ & & J_p \end{bmatrix} \otimes A \right) (T^{-1} \otimes I_{N_{vi}}) \\ &= (T \otimes I_{N_{vi}}) \begin{bmatrix} I_{n_1} \otimes M + \Delta t J_1 \otimes A & & \\ & \ddots & \\ & & I_{n_p} \otimes M + \Delta t J_p \otimes A \end{bmatrix} (T^{-1} \otimes I_{N_{vi}}), \end{aligned}$$

where  $n_j$  is the size of  $J_j$  (with  $\sum_{j=1}^p n_j = s$ ). Since the systems with the coefficient matrix  $T \otimes I_{N_{vi}}$  or  $T^{-1} \otimes I_{N_{vi}}$  can be solved directly and efficiently, we only need to consider the iterative solution of the systems associated with

$$I_{n_j} \otimes M + \Delta t J_j \otimes A, \quad j = 1, \dots, p. \quad (51)$$

For a Jordan block in the form (47), the above matrix becomes

$$I_{n_j} \otimes M + \Delta t J_j \otimes A = \begin{bmatrix} M + \mu_j \Delta t A & \Delta t A & & \\ & M + \mu_j \Delta t A & \ddots & \\ & & \ddots & \Delta t A \\ & & & M + \mu_j \Delta t A \end{bmatrix},$$

which can be solved by successively solving  $n_j$  systems with  $M + \mu_j \Delta t A$  (backward substitution).

On the other hand, for a Jordan block in the form (48), the matrix (51) becomes

$$I_{n_j} \otimes M + \Delta t J_j \otimes A = \begin{bmatrix} I_2 \otimes M + \Delta t C_j \otimes A & \Delta t I_2 \otimes A & & \\ & I_2 \otimes M + \Delta t C_j \otimes A & \ddots & \\ & & \ddots & \Delta t I_2 \otimes A \\ & & & I_2 \otimes M + \Delta t C_j \otimes A \end{bmatrix},$$

which, again, can be solved by successively solving  $n_j/2$  systems with  $I_2 \otimes M + \Delta t C_j \otimes A$ . Hence, the conditioning in this case is determined by the matrices in (49).

The analysis of Sect. 4 can be used to estimate the condition number of  $M + \mu_j \Delta t A$ . Estimate (50) for the second matrix in (49) is obtained similarly as for Theorem 5.1, except that we now have  $\sigma_{\max}(C_j) = \sigma_{\min}(C_j) = \sqrt{\alpha_j^2 + \beta_j^2}$ .  $\square$

Note that Remark 5.1 applies to the current situation as well.

## 6 Numerical examples

In the following examples we consider the IBVP (1) in 2d ( $d = 2$ ) with  $\Omega = (0, 1)^2$  and homogeneous Dirichlet boundary conditions. For the time integration we choose Radau5 (a Radau IIA method of order 5) with fixed time steps  $\Delta t = 10^{-1}$ ,  $10^{-3}$ ,  $10^{-5}$ , and  $\Delta t = N^{-\frac{1}{d}} \sim h$ , where  $h$  is the average mesh size. The singular values of the Radau5 coefficient matrix are  $\sigma_{\max} \approx 0.80$  and  $\sigma_{\min} \approx 0.09$  (see Table 1).

*Example 6.1.* To compare the condition number  $\kappa(I_s \otimes M + \Delta t \Gamma \otimes A)$  of Radau5 with the condition number  $\kappa(M + \Delta t A)$  of the implicit Euler method we consider two diffusion matrices: isotropic

$$\mathbb{D} = I \quad (\text{Laplace operator})$$

and anisotropic

$$\mathbb{D}(x, y) = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix} \begin{bmatrix} 10 & 0 \\ 0 & 0.1 \end{bmatrix} \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix}, \quad \theta = \pi \sin x \cos y, \quad (52)$$

with quasi-uniform (Fig. 1a) and  $\mathbb{D}^{-1}$ -uniform (Fig. 1b) meshes obtained using the mesh generator *bamg* [10]. A quasi-uniform mesh can be also seen as  $\mathbb{D}^{-1}$ -uniform for  $\mathbb{D} = I$ .

We expect the conditioning of Radau5 to be similar to that of the implicit Euler (see (43) and Remark 5.1). This is verified in Fig. 2 (left column), which shows that the conditioning for the Radau5 method has the same general behaviour as that for the implicit Euler scheme. With increasing number of mesh elements, the order of conditioning is  $\mathcal{O}(N^{\frac{2}{d}})$  for a fixed  $\Delta t$  and becomes  $\mathcal{O}(N^{\frac{1}{d}})$  for  $\Delta t \sim h \sim N^{-\frac{1}{d}}$  for both methods, which is in perfect agreement with the theoretical prediction in Remark 4.2.

Adapting towards the diffusion matrix improves the conditioning of the stiffness matrix [17]. This effect is observed in our test case as well: for anisotropic diffusion, the conditioning for diffusion-adapted meshes (Fig. 2c) is noticeably smaller than that with quasi-uniform meshes (Fig. 2b), especially as  $N$  is getting larger.

To compare the exact values for the conditioning of Radau5 with the bound (43) in Theorem 5.1, we use the exact eigenvalues for the matrices  $P^{-1/2} M P^{-1/2}$  and  $P^{-1/2} A P^{-1/2}$ . Figure 2 (right column) shows that the bound (43) is greater than the exact condition number but exhibits essentially the same behaviour as  $N$  increases.



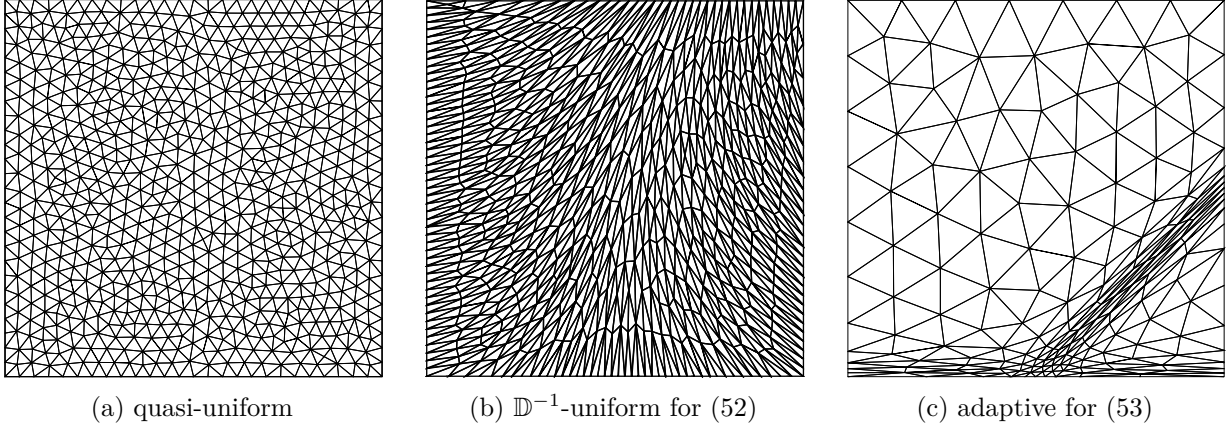


Figure 1: Mesh examples.

*Example 6.2.* To illustrate the effect of the diagonal scaling, we consider  $\mathbb{D} = I$  (Laplace operator),  $P = M_D + \Delta t A_D$ ,  $M_D$  and  $M_{lump}$ , and adaptive anisotropic meshes for the interpolation of the function

$$u(x, y) = \tanh(60y) - \tanh(60(x - y - 0.5)), \quad (x, y) \in (0, 1)^2, \quad (53)$$

generated with *bamg* [10] using an anisotropic adaptive metric [12]. This function simulates the interaction between a boundary layer along the  $x$ -axis and a shock wave along the line  $x = y + 0.5$  (see Fig. 1c for a mesh example).

Consistently with Remark 4.1, Fig. 3a shows that the diagonal scaling with  $P = M_D + \Delta t A_D$  (Euler) and  $I_s \otimes M_D + \Delta t \Gamma \otimes A_D$  (Radau5) reduces the effects caused by the mesh nonuniformity: the conditioning of both Euler and Radau5 methods is reduced by a factor ranging from 3 ( $\Delta t = 10^{-1}$ ) to 8 ( $\Delta t = 10^{-3}, 10^{-5}$ ).

For the scaling with  $P = M_D$ , Fig. 3b shows that the conditioning is getting worse for large  $\Delta t$  ( $\Delta t = 10^{-1}$ ) but is improving for small  $\Delta t$  ( $\Delta t = 10^{-3}, \Delta t = 10^{-5}$ ). To explain this, we recall that the diagonal entries of  $M$  are proportional to the corresponding patch volumes (see Lemma 3.1). Hence, scaling with  $M_D$  improves the conditioning issues caused by the mesh *volume-nonuniformity* (the same applies to  $M$  and  $M_{lump}$ ). On the other hand,  $\max_j A_{jj} \leq \lambda_{\max}(A) \leq (d+1) \max_j A_{jj}$  [17, Lemma 4.1] and Lemma 3.6 imply that the largest eigenvalue of  $A$  depends on the element shape (see also [17, section 4.1] and [19, section 3]). Thus,  $\kappa(M^{-1}A)$  is not necessarily smaller than  $\kappa(A)$  since rescaling of  $A$  with respect to the patch volumes does not necessarily improve the conditioning issues caused by the shape. Therefore, in general, the overall effect depends on the magnitude of  $\Delta t$ :  $\kappa(I + \Delta t M^{-1}A)$  is not necessarily better than  $\kappa(M + \Delta t A)$  for large  $\Delta t$ , but we can expect that  $\kappa(I + \Delta t M^{-1}A) < \kappa(M + \Delta t A)$  for small  $\Delta t$ .

## 7 Conclusions

The conditioning of implicit RK methods for the FE equations is comparable to the implicit Euler integration (Theorem 5.1). For the successive solution procedure, the conditioning of the system matrix  $I_s \otimes M + \Delta t \Gamma \otimes A$  of an implicit RK integration is determined by two types of smaller matrices (Theorem 5.2)

$$M + \mu_j \Delta t A \quad \text{and} \quad I_2 \otimes M + \Delta t \begin{bmatrix} \alpha_j & \beta_j \\ -\beta_j & \alpha_j \end{bmatrix} \otimes A.$$

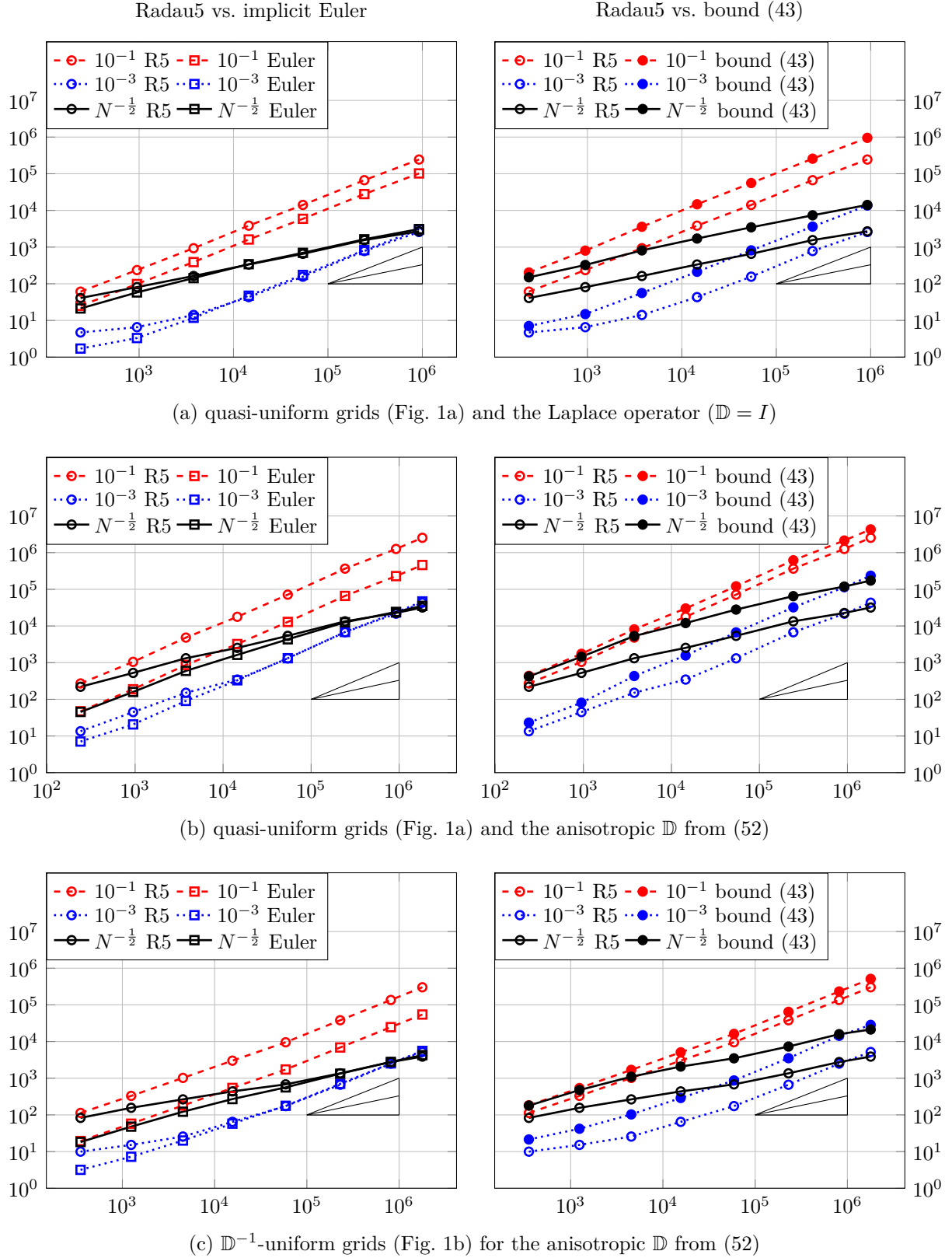


Figure 2: Conditioning for Radau5 (R5) and implicit Euler (left column) and a comparison of Radau5 conditioning with the bound (43) (right column) as functions of  $N$  (Example 6.1).

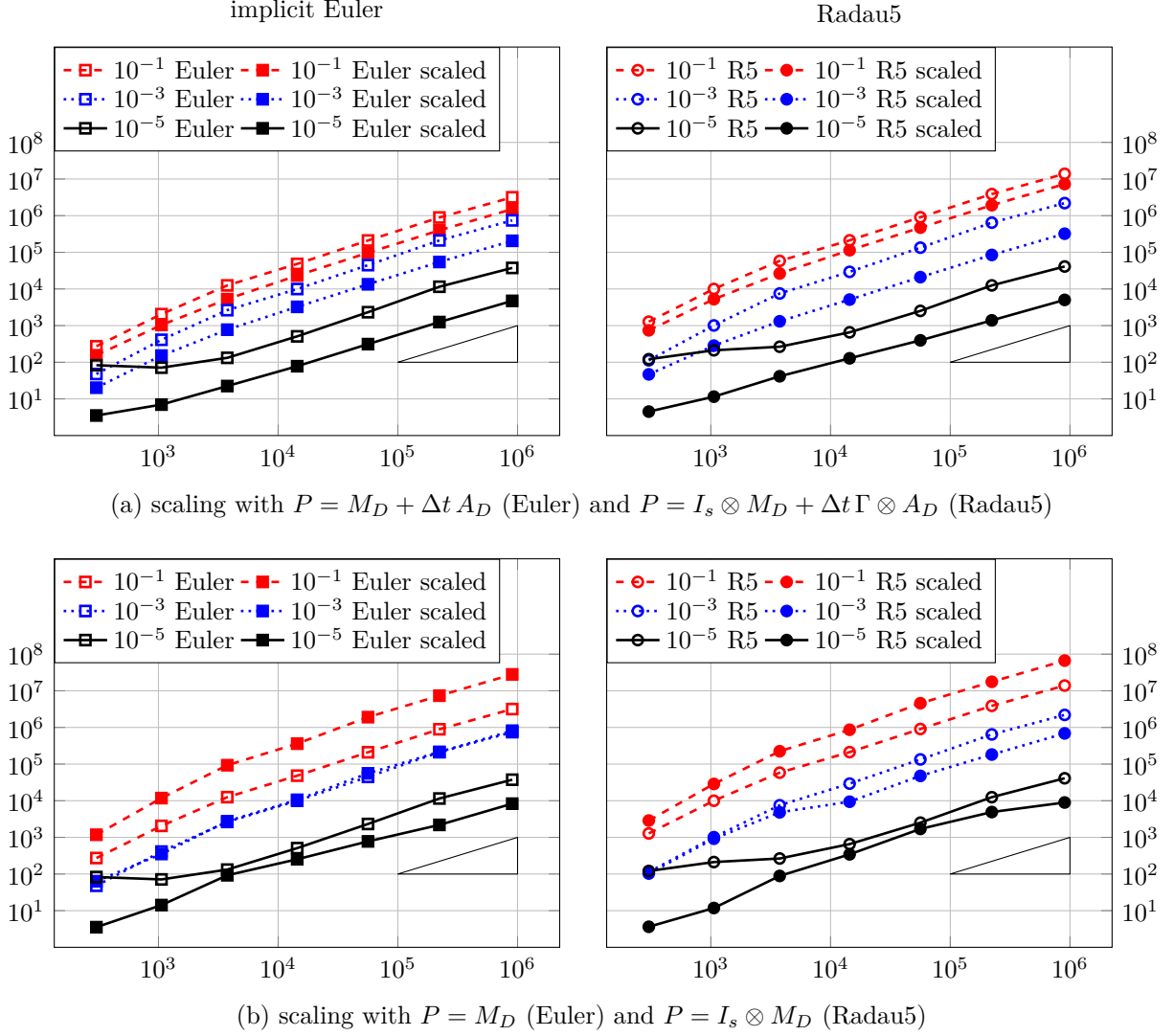


Figure 3: Conditioning of implicit Euler and Radau5 as function of  $N$  before and after a diagonal scaling for the Laplace operator and adaptive anisotropic meshes (Example 6.2).

The first matrix is similar to the implicit Euler method and corresponds to the real eigenvalues  $\mu_j$  of the RK matrix  $\Gamma$  while the second corresponds to the complex eigenvalues  $\alpha_j \pm i\beta_j$  of  $\Gamma$ .

There are three mesh-dependent factors that affect the conditioning of general implicit RK methods: the number of mesh elements (average mesh size), the mesh nonuniformity (in the Euclidean metric), and the mesh nonuniformity with respect to the inverse of the mass matrix. Preconditioning by the diagonal part of the system matrix itself or the mass matrix or by the lumped mass matrix reduces the effects of the mesh nonuniformity in the Euclidean metric (Theorems 4.2 and 4.3). These results are consistent with previous studies for the boundary value problems [1, 2, 8, 17] and explicit integration of linear diffusion problems [14, 15].

## References

- [1] M. Ainsworth, W. McLean, and T. Tran. The conditioning of boundary element equations on locally refined meshes and preconditioning by diagonal scaling. *SIAM J. Numer. Anal.*, 36(6):1901–1932 (electronic), 1999.
- [2] R. E. Bank and L. R. Scott. On the conditioning of finite element equations with highly refined meshes. *SIAM J. Numer. Anal.*, 26(6):1383–1394, 1989.
- [3] T. A. Bickart. An efficient solution process for implicit Runge-Kutta methods. *SIAM J. Numer. Anal.*, 14(6):1022–1027, 1977.
- [4] J. C. Butcher. On the implementation of implicit Runge-Kutta methods. *Nordisk Tidskr. Informationsbehandling (BIT)*, 16(3):237–240, 1976.
- [5] Q. Du, D. Wang, and L. Zhu. On mesh geometry and stiffness matrix conditioning for general finite element spaces. *SIAM J. Numer. Anal.*, 47(2):1421–1444, 2009.
- [6] G. J. Fix. Eigenvalue approximation by the finite element method. *Advances in Math.*, 10:300–316, 1973.
- [7] I. Fried. Bounds on the spectral and maximum norms of the finite element stiffness, flexibility and mass matrices. *Internat. J. Solids and Structures*, 9:1013–1034, 1973.
- [8] I. G. Graham and W. McLean. Anisotropic mesh refinement: the conditioning of Galerkin boundary element matrices and simple preconditioners. *SIAM J. Numer. Anal.*, 44(4):1487–1513 (electronic), 2006.
- [9] E. Hairer and G. Wanner. *Solving ordinary differential equations. II*, volume 14 of *Springer Series in Computational Mathematics*. Springer-Verlag, Berlin, second edition, 1996. Stiff and differential-algebraic problems.
- [10] F. Hecht. Bamg: Bidimensional Anisotropic Mesh Generator. <https://www.ljll.math.upmc.fr/hecht/ftp/bamg>.
- [11] N. J. Higham. *Accuracy and stability of numerical algorithms*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 1996.
- [12] W. Huang. Mathematical principles of anisotropic mesh adaptation. *Commun. Comput. Phys.*, 1(2):276–310, 2006.
- [13] W. Huang and L. Kamenski. On the mesh nonsingularity of the moving mesh PDE method. *Math. Comp.*, forthcoming. arXiv:1512.04971.
- [14] W. Huang, L. Kamenski, and J. Lang. Stability of explicit Runge-Kutta methods for high order finite element approximation of linear parabolic equations. In A. Abdulle, S. Deparis, D. Kressner, F. Nobile, and M. Picasso, editors, *Numerical Mathematics and Advanced Applications — ENUMATH 2013*, volume 103 of *Lecture Notes in Computational Science and Engineering*, pages 165–173. Springer International Publishing, 2015.
- [15] W. Huang, L. Kamenski, and J. Lang. Stability of explicit one-step methods for P1-finite element approximation of linear diffusion equations on anisotropic meshes. *SIAM J. Numer. Anal.*, 54(3):1612–1634, 2016.

- [16] L. Kamenski and W. Huang. A study on the conditioning of finite element equations with arbitrary anisotropic meshes via a density function approach. *J. Math. Study*, 47(2):151–172, 2014. arXiv:1302.6868.
- [17] L. Kamenski, W. Huang, and H. Xu. Conditioning of finite element equations with arbitrary anisotropic meshes. *Math. Comp.*, 83:2187–2211, 2014.
- [18] A. J. Laub. *Matrix analysis for scientists & engineers*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2005.
- [19] J. R. Shewchuk. What is a good linear element? Interpolation, conditioning, and quality measures. In *Proceedings of the 11th International Meshing Roundtable*, pages 115–126, Sandia National Laboratories, 2002.
- [20] L. Zhu and Q. Du. Mesh-dependent stability for finite element approximations of parabolic equations with mass lumping. *J. Comput. Appl. Math.*, 236(5):801–811, 2011.
- [21] L. Zhu and Q. Du. Mesh dependent stability and condition number estimates for finite element approximations of parabolic problems. *Math. Comp.*, 83(285):37–64, 2014.